

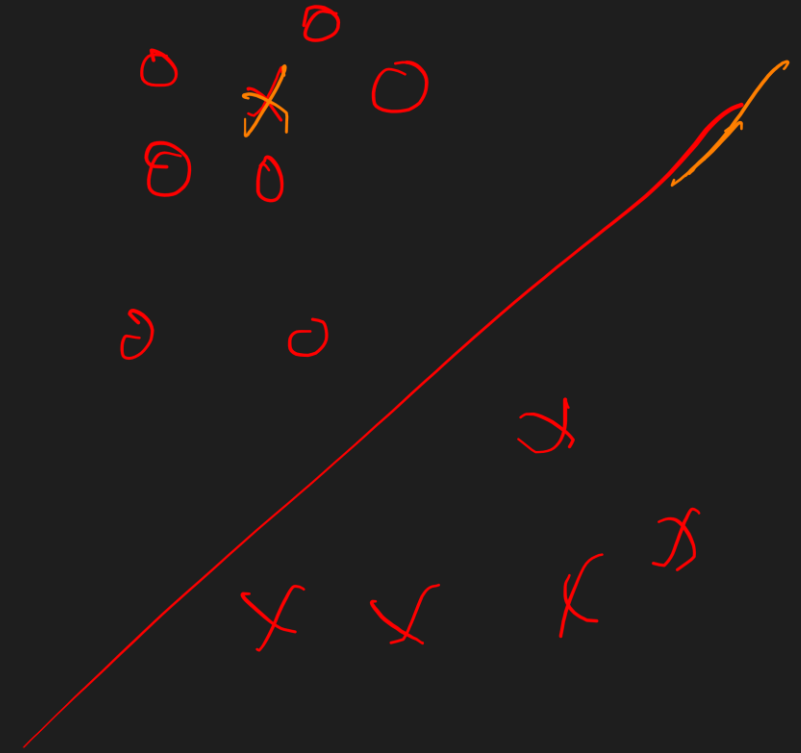
INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

LECTURE 14

HASSAN ASHTIANI

LINEAR CLASSIFIER

- $\hat{y}_w(x) = \text{sign}(W^T x) = \mathbf{1}_{W^T x \geq 0}$
- HOW TO FIND W GIVEN DATA SET Z ?
 - $\hat{W} = \arg\text{MIN}_w L_Z^{0-1}(\hat{y}_w)$
 - HOW TO MINIMIZE?



HARDNESS OF LINEAR CLASSIFICATION

- IN GENERAL, FINDING THE LINEAR SEPARATOR WITH MINIMUM CLASSIFICATION ERROR IS NP-HARD (WITH RESPECT TO d)

- UNLESS....

- DATA IS LINEARLY SEPARABLE! ✓

- ..OR OPTIMIZE A SURROGATE LOSS FUNCTION INSTEAD! ✓

- CAN'T WE JUST GO INTO HIGHER DIMENSIONS TO MAKE THE DATA LINEARLY SEPARABLE?

$$L^{d-1} (y, \hat{y}) = \begin{cases} 1 & \text{a.w} \\ 0 & \hat{y} = y \end{cases}$$

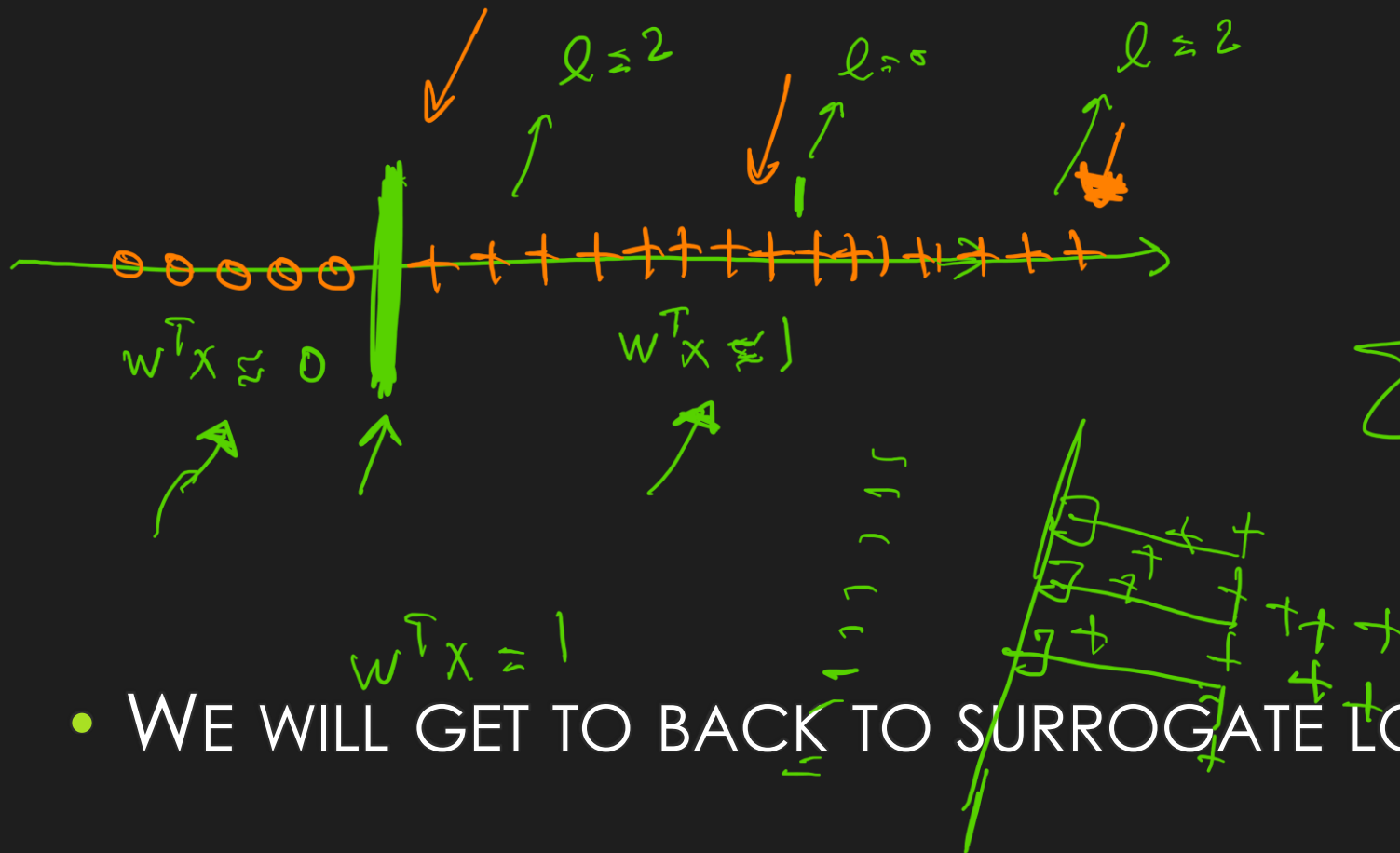
SURROGATE LOSS FUNCTIONS

- SO FAR, WE ASSUMED THE CLASSIFIER RETURNS A DISCRETE VALUE
 - E.G., $\hat{y}_w = \text{sign}(W^T x) \in \{0,1\}$
- WHAT IF THE CLASSIFIER'S OUTPUT IS CONTINUOUS
 - E.G., $\hat{y}_w = W^T x$
 - \hat{y} WILL CAPTURE THE "CONFIDENCE" OF THE CLASSIFIER TOO
- OTHER (CONTINUOUS) LOSS FUNCTIONS
 - MARGIN LOSS, CROSS-ENTROPY/NEGATIVE-LOG-LIKELIHOOD LOSS, ...
- POTENTIAL BENEFITS?
 - EASIER TO OPTIMIZE
 - MITIGATE OVERFITTING

$\text{sign}(w^T x)$

SQUARED LOSS FOR CLASSIFICATION?

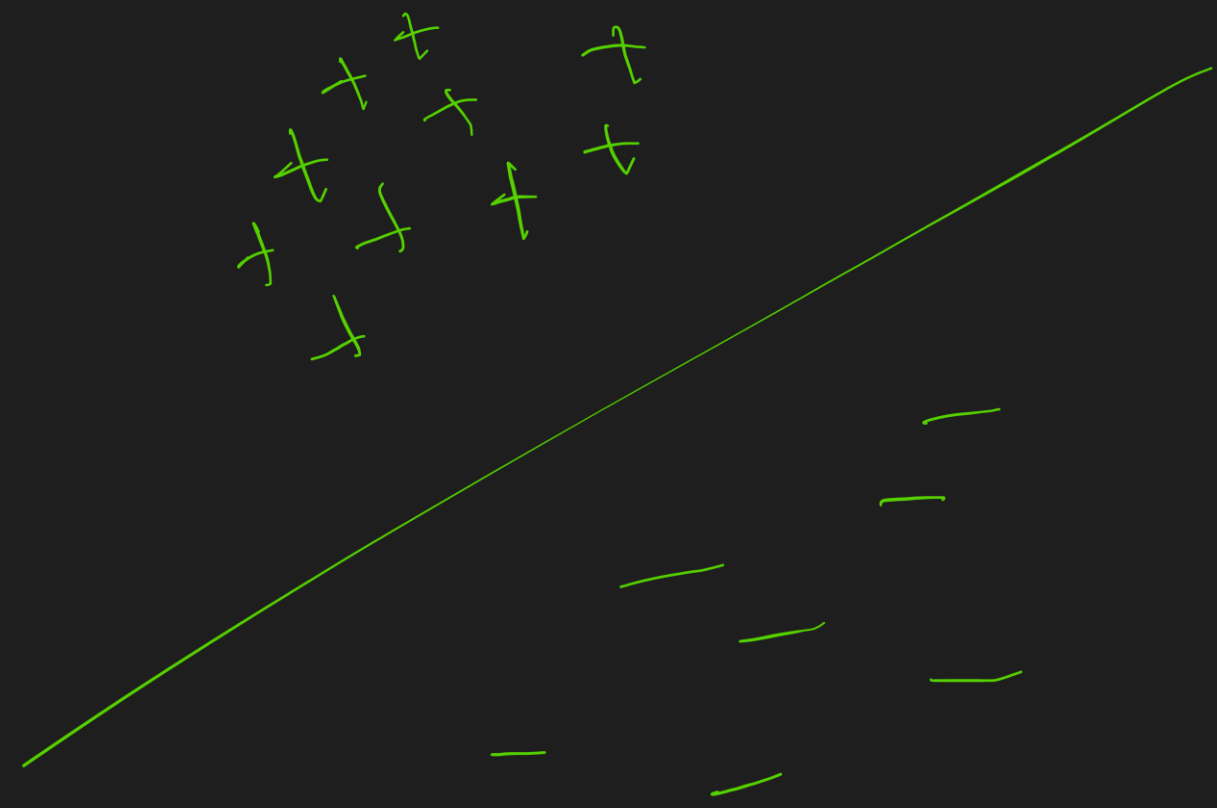
$x_1 = 10, y_1 = 1$
 $x_2 = 11, y_2 = 1$
 $x_3 = 8, y_3 = 0$



$$\sum (\hat{y}^i - y^i)^2 = \sum (w^T x^i - y^i)^2$$

- WE WILL GET TO BACK TO SURROGATE LOSS FUNCTIONS

LINEARLY SEPARABLE DATA



$$y^i \in \{+1, -1\}$$

LINEARLY SEPARABLE DATA

- A BINARY CLASSIFICATION DATA SET $Z = \{(x^i, y^i)\}_{i=1}^n$ IS LINEARLY SEPARABLE IF
 - THERE EXISTS W^* SUCH THAT
 - FOR EVERY $i \in [n]$ WE HAVE $\text{sgn}(\langle x^i, W^* \rangle) = y^i$
 - OR EQUIVALENTLY, FOR EVERY $i \in [n]$ WE HAVE $(W^{*T} x^i) y^i > 0$
- IN OTHER WORDS, THE CLASSIFICATION ERROR ON Z IS 0
- CAN WE FIND W^* EFFICIENTLY FOR LINEARLY SEPARABLE DATA?

LINEAR PROGRAMMING

- STANDARD LP PROBLEM:

$$\begin{aligned} & \downarrow \\ & \max_{w \in \mathbb{R}^d} \langle u, w \rangle = \\ & \text{s.t. } Aw \geq v \end{aligned}$$

~~Given~~ A, u, v are given.

- LP PROBLEMS CAN BE SOLVED EFFICIENTLY!

$$\rightarrow u, w \in \mathbb{R}^d$$

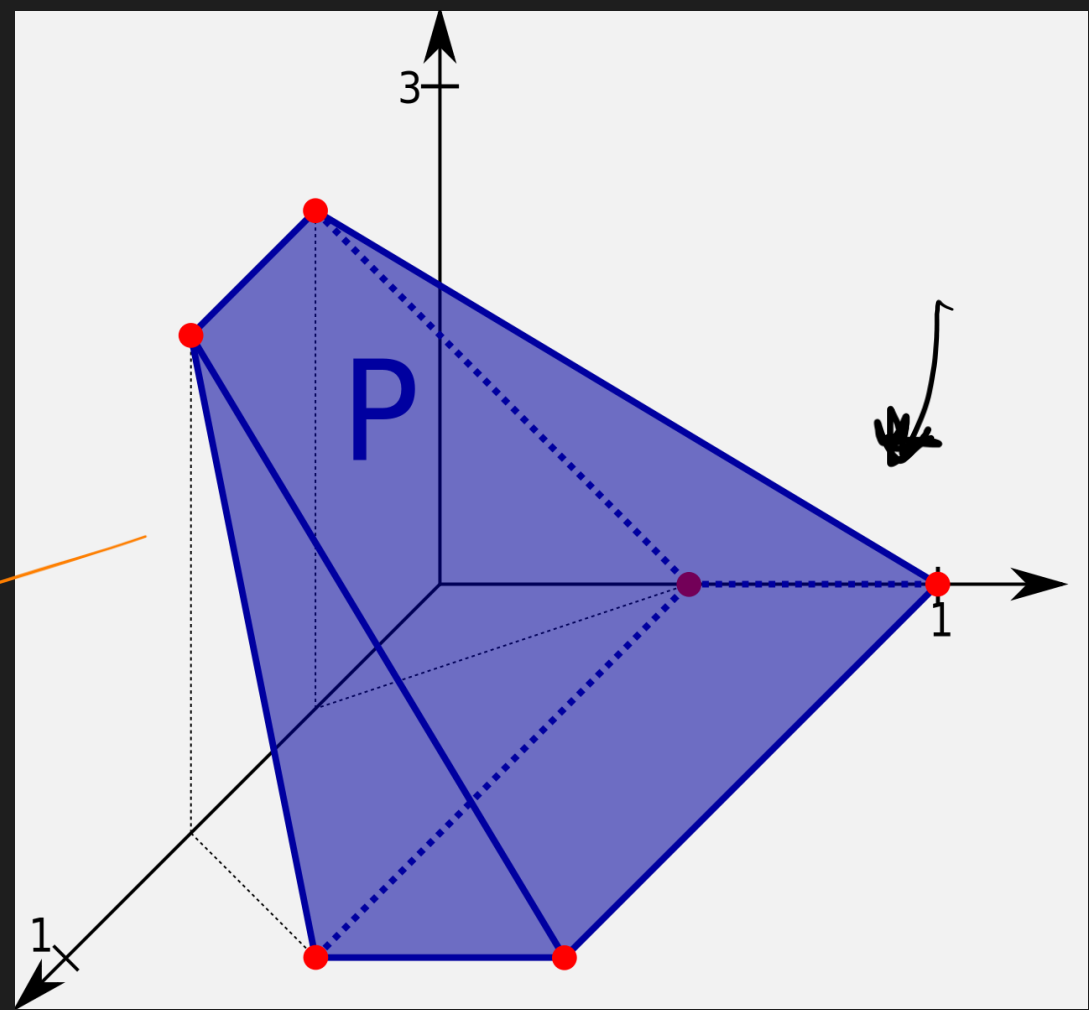
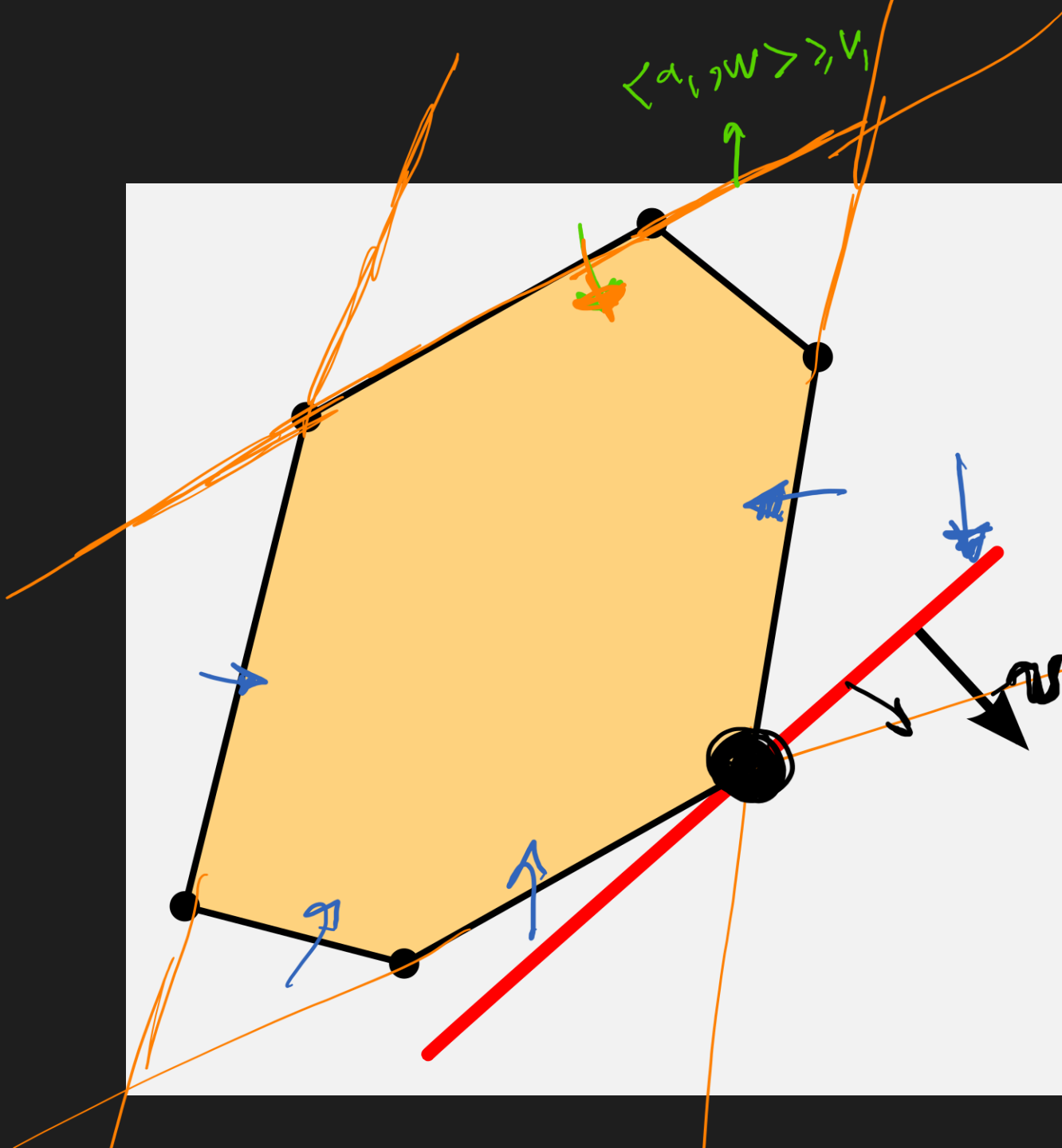
$A_{k \times d}$ matrix, $v \in \mathbb{R}^k$

$$A = \begin{bmatrix} a_1^T \\ \vdots \\ a_k^T \end{bmatrix}_{k \times d}$$

$$\sum_{i=1}^d u_i w_i$$

k linear constraints

$$\begin{bmatrix} \langle a_1, w \rangle \geq v_1 \\ \vdots \\ \langle a_k, w \rangle \geq v_k \end{bmatrix}$$



LP FOR CLASSIFICATION

$$(W^{*T} x^i) y^i = \delta \cdot \delta$$

- DATA IS LINEARLY SEPARABLE SO

→ • $\exists W^*$ s.t. $\forall i \in [n], (W^{*T} x^i) y^i > 0$

- SO,

• $\exists W^*, \gamma > 0$ s.t. $\forall i \in [n], (W^{*T} x^i) y^i \geq \gamma$

- SO,

• $\exists W^*$, s.t. $\forall i \in [n], (W^{*T} x^i) y^i \geq 1$

$$W_{new}^* \approx W_{old}^* \times \frac{1}{\delta}$$

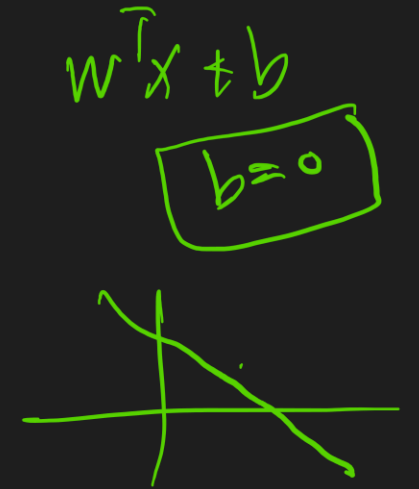
LP FOR LINEAR CLASSIFICATION

- DEFINE $A = [x_j^i y^i]_{n \times d}$
- THEN FINDING THE OPTIMAL W IS EQUIVALENT TO

$$\max_{w \in \mathbb{R}^d} \langle \vec{0}, w \rangle$$

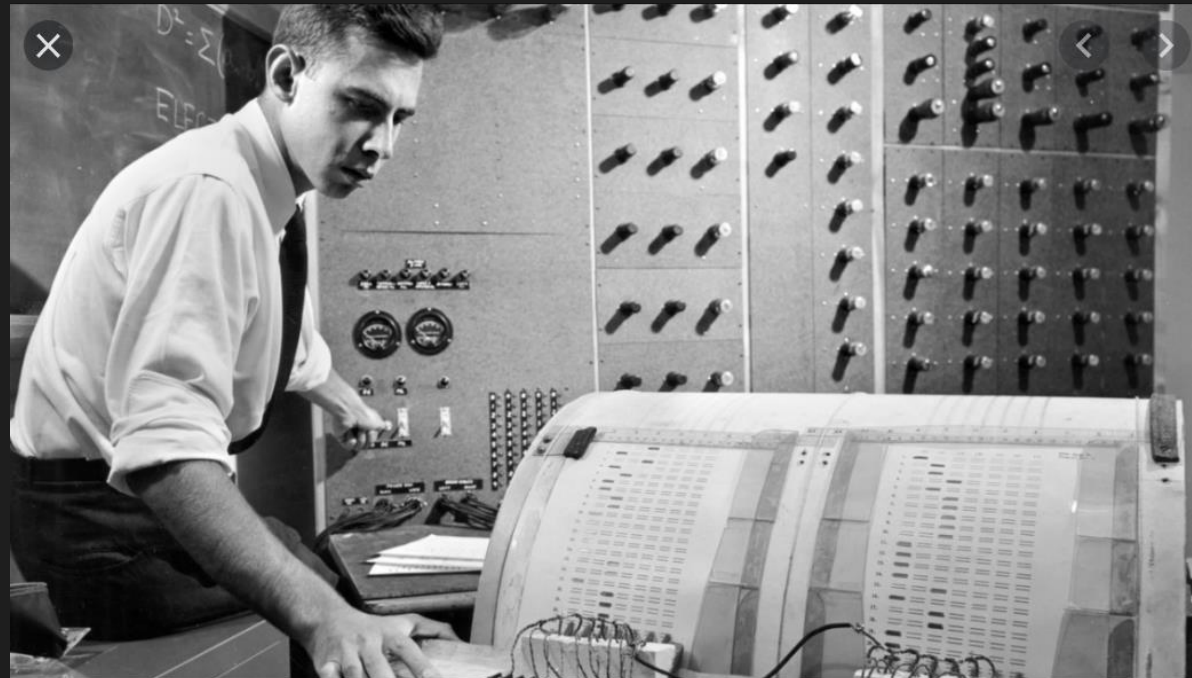
$$\text{s. t. } \underline{Aw} \geq \underline{\vec{1}}$$

- WE CAN USE OFF-THE-SHELF LP SOLVERS.
- WHAT IF THE BEST W DOES NOT GO THROUGH THE ORIGIN? (IT HAS A BIAS OR INTERCEPT)? \rightarrow



APPROACH 2: PERCEPTRON

- PROPOSED IN 50'S BY ROSENBLATT
- PREDECESSOR OF NEURAL NETWORKS
 - MULTI-LAYER PERCEPTRON!



ROSENBLATT'S PERCEPTRON

Batch Perceptron

input: A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

initialize: $\mathbf{w}^{(1)} = (0, \dots, 0)$

→ **for** $t = 1, 2, \dots$

if $(\exists i \text{ s.t. } y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle \leq 0)$ then

→ $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \mathbf{y}_i \mathbf{x}_i$

else

output $\mathbf{w}^{(t)}$

- IN EACH UPDATE, W BECOMES “MORE CORRECT” ON x^i
- [HTTPS://PHIRESKY.GITHUB.IO/KOGSYS-DEMOS/NEURAL-NETWORK-DEMO/?PRESET=ROSENBLATT+PERCEPTRON](https://phiresky.github.io/kogsys-demos/neural-network-demo/?preset=rosenblatt+perceptron)

$$y^i \in \{+1, -1\}$$

THE GREEDY UPDATE

- IN EACH UPDATE, W BECOMES "MORE CORRECT" ON x^i :

$$\underline{W_{new}^T x^i y^i} = \langle \underline{w_{old} + y^i x^i}, x^i \rangle y^i$$

$$= w_{old}^T x^i y^i + \|x^i\|_2^2 y^i y^i$$

- WHAT ABOUT OTHER x^j 'S?

$$\gg \underline{w_{old}^T x^i y^i}$$